**Credit Risk Modeling Proposal**
Implementing a machine learning (ML)-based credit risk system will enhance Citi's loan management by automating risk assessment, reducing defaults, and enabling data-driven decisions. This system will improve accuracy over traditional methods, optimize interest rates, and ensure compliance with regulatory standards.

---

## Data Requirements

**Input Variables:**

- **Customer Data**:
    - Credit score, age, employment status, income, debt-to-income (DTI) ratio.
    - Payment history, existing liabilities, collateral value.
- **Loan Details**:
    - Loan amount, term, purpose (e.g., mortgage, personal).
- **Macroeconomic Indicators**:
    - Unemployment rate, inflation, interest rates.
- **Behavioral Data**:
    - Transaction patterns, savings habits.
- **External Data**:
    - Credit bureau reports (e.g., Experian), public records (bankruptcies).

**Data Sources**:

- Internal databases (application forms, transaction history).
- External APIs (credit bureaus, economic datasets).

---

## Data Outputs

- **Risk Probability Score**:
    - Probability of default (e.g., 0–100%).
- **Risk Classification**:
    - Labels: *Low Risk*, *Medium Risk*, *High Risk*.
- **Recommended Actions**:
    - Loan approval/rejection, adjusted interest rates, collateral requirements.
- **Explainability Reports**:
    - SHAP values or LIME outputs to justify decisions (critical for regulatory compliance).

---

## Architecture

**Model Selection**:

- **Gradient Boosting Machines (XGBoost/LightGBM)**:
    - Handles non-linear relationships, robust to missing data.
    - Provides feature importance scores.
- **Hybrid Approach**:

    - Combine ML with logistic regression for interpretability.
  - **Deep Learning (Optional)**:
    - Neural networks for unstructured data (e.g., text-based employment history).

**Pipeline**:

1. **Data Ingestion**: Batch/real-time data collection.
2. **Preprocessing**: Handle missing values, normalize features.
3. **Feature Engineering**: Derive metrics like DTI, payment consistency.
4. **Model Training**: Cross-validation to prevent overfitting.
5. **API Deployment**: Integrate with Citi's loan management system for real-time scoring.
6. **Monitoring**: Track accuracy, fairness, and drift over time.

**Tech Stack**:

- Python (scikit-learn, XGBoost), TensorFlow/Keras (for DL).
- AWS/GCP for scalable compute.
- Airflow for pipeline orchestration.

---

# Risks and Challenges

1. **Data Quality**: Missing/inaccurate historical data may skew predictions.
2. **Bias & Fairness**: Models might inherit biases from past discriminatory practices.
   - Mitigation: Regular fairness audits, bias-correction algorithms.
3. **Regulatory Compliance**: GDPR, ECOA, and "right to explanation" requirements.
4. **Model Drift**: Economic shifts (e.g., recessions) degrade performance.
   - Mitigation: Retrain models quarterly with updated data.
5. **Interpretability**: Balancing accuracy with explainability for regulators and customers.
6. **Integration Costs**: Compatibility with legacy systems.

---

# Conclusion

This proposal outlines a scalable, compliant credit risk system that leverages ML to minimize defaults while maintaining transparency. Next steps include a pilot program with historical data validation and stakeholder training.

**Estimated Impact**:

- 20–30% reduction in default rates.
- 15% faster loan approval turnaround.
- Improved customer trust through explainable decisions.

---